

Questimator: Generating Knowledge Assessments for Arbitrary Topics

Qi Guo, Chinmay Kulkarni, Aniket Kittur, Jeffrey P. Bigham and Emma Brunskill
School of Computer Science, Carnegie Mellon University
qigu@andrew.cmu.edu, {chinmayk, nkittur, jbigham, ebrunskill}@cs.cmu.edu

Abstract

Formative assessments allow learners to quickly identify knowledge gaps. In traditional educational settings, expert instructors can create assessments, but in informal learning environment, it is difficult for novice learners to self assess because they don't know what they don't know. This paper introduces Questimator, an automated system that generates multiple-choice assessment questions for any topic contained within Wikipedia. Given a topic, Questimator traverses the Wikipedia graph to find and rank related topics, and uses article text to form questions, answers and distractor options. In a study with 833 participants from Mechanical Turk, we found that participants' scores on Questimator-generated quizzes correlated well with their scores on existing online quizzes on topics ranging from philosophy to economics. Also Questimator generates questions with comparable discriminatory power as existing online quizzes. Our results suggest Questimator may be useful for assessing learning in topics for which there is not an existing quiz.

1 Introduction

An increasing number of learners are opportunistic and self-driven [Brandt *et al.*, 2009; Sadowski *et al.*, 2015; Cai *et al.*, 2015]. Online learning opportunities such as MOOCs attract diverse students with a range of prior knowledge and experience [Kulkarni *et al.*, 2015]. While online resources are a boon to self-directed learning, knowledge gaps cause learners to struggle with new material, lose motivation, and even avoid future subjects due to lowered self-perceived efficacy [Horwitz and Cope, 1986].

Assessments can help learners diagnose what they know, which may be particularly helpful in informal learning domains. Unfortunately informal learning domains are precisely those that often lack formal assessments. Relying on learners' self-diagnosis is tricky because they don't know what they don't know [Caputo and Dunning, 2005].

Traditionally, formative assessment of factual knowledge has relied on experts, such as teachers or textbook authors. However, what such expert resources provide in quality, they

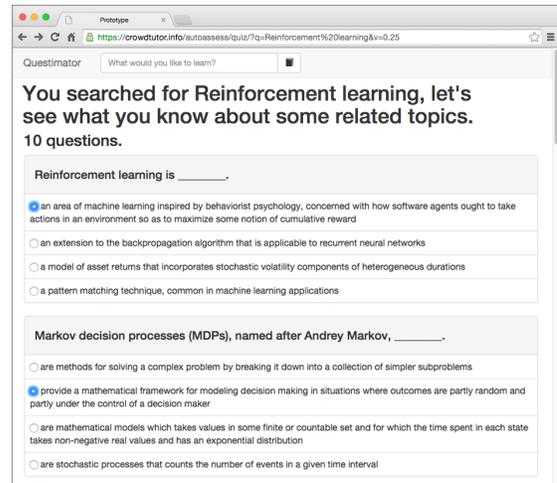


Figure 1: Questimator automatically produces quizzes of multiple-choice questions to assess knowledge for arbitrary topics contained within Wikipedia. (The selected answers in the figure are correct answers.)

lack in breadth [Elwood and Klenowski, 2002]. Because creating formative assessments of factual knowledge is time-consuming, such assessments are limited to a small set of common topics of interest. Furthermore, such assessments are often missing for new or esoteric topics.

This paper introduces Questimator (described in Section 3), which, to our knowledge, is the first system that automatically generates multiple-choice question (MCQ) quizzes for general topics from large networked corpora like Wikipedia, thus significantly extending the scope of automatically generated questions. For example, Figure 1 shows Questimator's automatically generated assessment for the topic of "Reinforcement learning".

A learner uses Questimator by inputting a topic they would like to learn about, *e.g.* "reinforcement learning" or "France". To generate a coherent quiz, we need questions covering different perspectives of that topic. Unfortunately, previous work (Section 2) mostly focused on single MCQ generation, and did not provide a principled quiz generation approach. We propose a method (*contribution 1*) to find topics that are related to the input topic (Section 3.1) and generate questions

for those related topics. We utilize the Wikipedia page-to-page linkage graph and word embedding in the process.

Generating a single multiple-choice question (Section 3.2) involves creating three pieces of output: (i) a question stem, (ii) a correct answer, and (iii) a set of distractors (usually, at least 3). A question stem (Section 3.3) is the part where the item to ask for is stated, for example “Reinforcement learning is _____.” The examinee then chooses from correct answer and distractors. We propose a novel method (*contribution 2*) to find distractors (Section 3.4), combining Wikipedia’s categorical structure (treated as a bipartite graph between the set of categorical labels and the set of articles), word embedding and sentence embedding.

By strengthening related topic search and distractor generation with Wikipedia’s graph structure and semantic embedding, Questimator advances question generation for general topics. Through Questimator, we demonstrate the feasibility of leveraging existing corpora of knowledge such as Wikipedia to create interactive testing materials.

2 Related Work

2.1 Automatic MCQ Generation

Questimator is most closely related to prior systems for generating *textual* multiple-choice questions, as creating questions for math or logic involves different techniques. Most prior methods for generating text-based questions start with an input article. The question and correct answer are derived from this article (often from a single sentence). And they mostly aimed at language learning [Hoshino and Nakagawa, 2005; Brown *et al.*, 2005; Mitkov *et al.*, 2006; Lin *et al.*, 2007; Smith *et al.*, 2010; Gates *et al.*, 2011; Mostow and Jang, 2012] or a single subject/textbook [Wang *et al.*, 2008; Agarwal and Mannem, 2011; Bhatia *et al.*, 2013].

These prior approaches to generating MCQ quizzes are limited to using single document or a fixed ontology alone to generate questions and distractors. Generating questions from a single document often scopes questions too narrowly to assess a learner’s understanding of a general topic, which typically spans related documents. Related topic selection is addressed in Section 3.1.

Also, while previous approaches work well for language learning or single subject, they suffer from limited choices of distractors when directly applied to general topics. Some systems pick distractors (having same POS tag as the correct answer) within the same article/textbook, based on term frequency etc. [Hoshino and Nakagawa, 2005; Agarwal and Mannem, 2011]. Using terms from the article as distractors can be problematic, because good distractors might not be frequently mentioned in the article. For example, “Formula Two” is a good distractor for “Formula One”, but does not have a high frequency in the Wikipedia article for “Formula One”. Other systems use fixed ontologies or dictionaries, such as *WordNet* [Miller, 1995], to find synonyms or other related words to use as distractors [Mitkov *et al.*, 2006; Gates *et al.*, 2011].

Although *WordNet* has 117,000 synsets, it focuses more on a limited type of cognitive concepts, rather than general knowledge concepts. For example, it contains “baby”, and

groups it with “infant”, but it does not contain “Baby”, the Justin Bieber hit song. This limitation also applies to other ontologies and dictionaries. In comparison, there are more than 5,129,000 English Wikipedia articles (by April 2016) for a broad range, containing long-tail topics and contemporary topics. We use the abundant Wikipedia corpora along with its socially-annotated categorical information [Kittur *et al.*, 2009] to select distractor topics. So Questimator can work across the wide range of domains. In addition it integrates semantic embeddings when generating distractors.

2.2 Crowdsourcing Question Generation

Crowdsourcing has recently been explored as a way to generate questions. Christoforaki and Ipeirotis [2015] ask crowdworkers to directly generate the questions by mining existing Q&A sites. For example, Stack Overflow can be used as a scalable source of coding questions. However this subjects to extra economical cost, human selection of topics and additional time delay. Another challenge with these approaches is verifying that the questions are of high quality. Although Wikipedia is crowdsourced, popular articles are quite accurate and complete, and so the questions generated by Questimator are in a sense crowdsourced from source material subject to Wikipedia’s quality control.

Questimator is a new approach to multiple-choice question generation. Inspired by the recent trend to utilize large scale datasets, and enabled by the corresponding development in NLP, Questimator uses a connected set of documents rather than a single document to generate multiple choice questions for general topics.

3 Questimator System

Questimator is our system for automatically generating multiple choice questions from Wikipedia. To use Questimator, learners first enter a topic they would like to know about (Questimator matches the entered topic to one that appears in Wikipedia). Questimator then returns a set of multiple choice questions automatically generated from Wikipedia chosen to assess the learner’s knowledge of the provided topic (Figure 1 shows the first two questions generated for the topic “Reinforcement learning”). The number of questions to return is configurable (10 by default). To do this, Questimator:

1. generates a list of topics related to the input topic,
2. generates questions for each of the topics by
 - (a) generating question stems,
 - (b) generating and ranking distractors for each question stem.

3.1 Stage 1: Related topic identification

Related topics for questions can be drawn from prerequisites and subtopics. Given an input (main) topic X_0 , to generate candidate related topics $X_i (i = 1, \dots, n)$, Questimator leverages both the Wikipedia document for the main topic, and the larger Wikipedia structure to generate and rank related topics.

To rank a candidate related topic T , Questimator uses three measures of similarity:

- Term frequency of T in the article of X_0

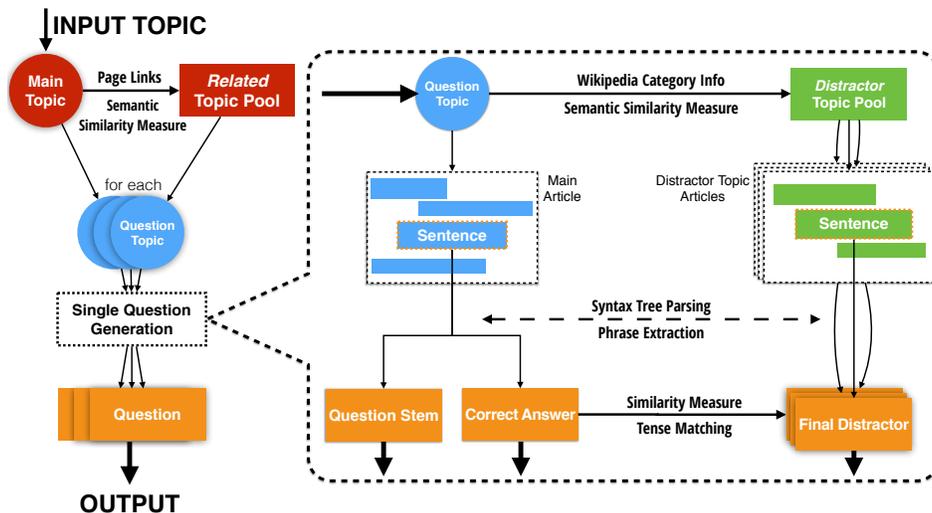


Figure 2: When learners enter a topic, Questimator extracts the wiki links on the page and finds related topics. A question is then generated for both the input topic and each of the related topics. To generate a question from a topic, categories from the Wikipedia page for the topic are used to generate similar topics. Sentences extracted from these pages are then turned into question stems and both correct answers and distractors. In the case of distractors, similarity and tense matching are used to pick the final distractors. The question stem, correct answer, and distractors form the final output question.

- backlink overlap between T and X_0
- embedded semantic similarity between T and X_0

We will elaborate on each of the three features in the following. Questimator combines these measures to produce a single relatedness metric. And Questimator then ranks topics based on that, and chooses the top n ($n \geq 9$) items as related topics to generate questions for the quiz.

Wikilink term frequency

Important concepts are likely to be repeatedly mentioned in the *main article* of X_0 , and so term frequency can be an important feature for identifying important subtopics. We extract all the links to other Wikipedia topics in the main article (called *pagelinks* or *wikilinks*), and count the term frequencies ($\text{TF}_{X_0}(T)$) of the wikilinks. We take simple confereces into account, like abbreviations used in the article etc.

The term frequencies alone do not always identify good related topics. For example, “(software) agent” are mentioned many times in the article for “Reinforcement learning”, but it is not an important related concept. Neither is “Real number” for “Invertible matrix”.

Backlink Overlap

The topics that co-occur with the main topic frequently are likely to be related. Similar backlinks, i.e. articles of which topics mention the topic, may imply topic relation. We find the backlinks of the candidate topics and the main topic, and compute the cosine similarity between their “bag-of-backlinks” vectors using the following formula:

$$\text{sim}_{\text{bl}}(T, X_0) = \frac{|\text{backlinks}(T) \cap \text{backlinks}(X_0)|}{\|\text{backlinks}(T)\|_2 \cdot \|\text{backlinks}(X_0)\|_2}$$

, where $|\cdot|$ is the cardinality of the set, X_0 is the main topic, T is any candidate topic. This allows us to identify topics that

are overlooked in the main article that can not be identified by term frequency. For example, “Convolutional neural network” for “Deep learning”, and “Affine transformation” for “Rigid transformation”.

Embedded semantic similarity of topics

We used a Word2Vec [Mikolov *et al.*, 2013] model trained on Wikipedia articles with topics (article titles) treated as independent entities. This produces a vector representation for each title and word. We then measure the similarity between topics T and X_0 by their vectors’ cosine similarity, $\text{sim}_{\text{word2vec}}(T, X_0) = \frac{\text{vec}(T) \cdot \text{vec}(X_0)}{\|\text{vec}(T)\|_2 \cdot \|\text{vec}(X_0)\|_2}$. Adding this feature to the relatedness measure, we might pick out “Back-propagation” for “Artificial neural network”, which does not have either high term frequency or backlink overlap, but is semantically related, as an important optimization method for neural networks.

Combine and Rank Relatedness

We combine the above three quantities, to produce a single scalar used for ranking. We normalize each of the three features independently by $\tilde{f} = \frac{f - \text{mean}(f)}{\sqrt{\text{var}(f)}}$. Then they are truncated to $[-1, 1]$, so a single feature would not play a too significant role in some edge cases. We sum the three features to create a combined relatedness measure,

$$R_{X_0}(T) = \widetilde{\text{TF}}_{X_0}(T) + \widetilde{\text{sim}}_{\text{bl}}(T, X_0) + \widetilde{\text{sim}}_{\text{word2vec}}(T, X_0)$$

. These features can be assigned different weights, here we simply set them to be equally weighted. Links are finally sorted by their relatedness values, and the top n (9 by default) are picked as related topics to generate questions. For example, the related topics picked for “Reinforcement learning”

are “Markov decision processes”, “Software agent”, “Dynamic programming”, “Temporal difference learning”, “Supervised learning”, “Machine learning”, “Optimal control theory”, “Gradient descent” and “Q-learning”.

3.2 Stage 2: Single question generation

Given a topic $X_i (i = 0, \dots, n)$ as input, Questimator generates a single MCQ as following:

- Question stem and correct answer are generated first from the same sentence, and
- distractors are generated separately from articles of similar topics.

3.3 Stage 2(a): Question stem generation

Questimator generates questions directly testing on concept explanation are frequently used in expert assessments, by asking about the verbal phrases after the topic noun phrases in the main clause. Questimator generates *gap-fill questions*, as the stems. To generate a question stem for a given topic, Questimator retrieves the Wikipedia article of the topic, and

- finds a set of sentences that each contain the stemmed tokens of the topic, and process them by their order in the article,
- parses a sentence into a PCFG syntax tree, and uses TGrep [Rohde, 2004] to match the syntax tree with certain syntactic phrase patterns,
- if matched, substitutes the matched phrase with a blank to generate the question stem, and uses the matched phrase as the correct answer.

Then Questimator processes the sentences by their order in the article, as the first few sentences where the main topic is mentioned usually contain an explanation of that topic. Note the second step is similar to Heilman and Smith [2010].

Additional types of question stem can be generated by simply inserting more TGrep patterns into Questimator. As a byproduct of limiting the question stem types, we avoid the problem of choosing the proper question type.

3.4 Stage 2(b): Distractor Generation

Different from previous approaches, we propose to generate distractors utilizing the Wikipedia categorical information, and word/sentence embedding methods. To generate distractors, Questimator

- finds topics in the same categories as the main topic,
- ranks them by their Word2Vec semantic similarities with the main topic, and picks the top few topics,
- extracts one distractor phrase for each by matching patterns on syntax tree (same as extracting the correct answer in Section 3.3),
- uses skip-thought vectors [Kiros *et al.*, 2015] to find distractor phrases most similar to the right answer phrase.

Pre-Extraction: Distractor Topic Selection

Unlike searching for *related topics* as question topics (Section 3.1), the criterion of similarity for selecting distractor topics is different. For example, “Camera Matrix” is a very related topic for “Camera Calibration”, by our former criteria for relatedness. But as a distractor, examinees can easily tell a *matrix*’s definition from a calibration *process*. Luckily,

semantic information in the definition phrase reveals this difference, so sentence embedding can help. In general, topics at the same level of the category hierarchy are preferred. For example, “Supervised learning” and “Unsupervised learning”, “Formula One” and “Formula Two”. So categorical structure is more informative than linkage graph for this purpose.

Questimator first finds topics sharing at least one category with the question topic to construct a candidate pool. The category information is from Wikipedia’s socially annotated (noisy) category hierarchy. Then Questimator ranks the distractor topic candidates by their Word2Vec similarity (same as in Section 3.1) with the question topic, and selects the top ones. For n_d ($n_d = 3$ by default) distractors, Questimator intermediately select $m \cdot n_d$ distractor topics ($m = 3$ by default) to generate the distractor phrases.

Post-Extraction: Distractor Phrase Ranking

After we have mn_d distractor phrases, we aim to pick the ones most difficult to distinguish from the correct answers. Questimator applies skip-thought vectors [Kiros *et al.*, 2015] to measure the similarity between the distractors and the correct answer. Like Word2Vec generating vectors for words, skip-thought generates vectors for sentences, with semantically similar sentences having similar vectors.

Questimator throws away phrases containing the stemmed tokens of the main topic string, because there is a high probability that this will reveal itself as a wrong answer. We match the distractors’ tenses with the correct answer. This prevents the examinees from identifying them as wrong answers by tense mismatch.

3.5 Questimator MCQ Example

Finally distractors together with the question stem and the correct answer are delivered as a whole MCQ. The following is an example of a whole MCQ (correct answer in italic):

A recurrent neural network (RNN) is _____.

- *a class of artificial neural network where connections between units form a directed cycle*
- an artificial neural network where connections between the units do not form a directed cycle
- a parallel computing paradigm similar to neural networks, with the difference that communication is allowed between neighbouring units only
- a type of artificial neural network in which an electrically adjustable resistance material is used to emulate the function of a neural synapse

One of our goals in developing Questimator is for it to generate MCQ quizzes at scale for topics that do not have quizzes available online. We maintain an updated corpus of quizzes generated from the most popular Wikipedia articles at <https://crowdtutor.info>. For many topics, such as “Wonders of the World” and “Oasis (band)”, no other existing quizzes are readily available. For topics like the ongoing “Syrian Civil War”, expert-generated quizzes, even if available, are likely out-of-date. Finally, for some topics, such as “Fascism” and “Constructivism”, we found that Questimator generated questions similar to those on existing quizzes.

4 Evaluation

The goal of Questimator is to generate good MCQ quizzes for knowledge evaluation for arbitrary topics.

To preliminarily evaluate Questimator questions' quality, we asked 4 TAs to inspect 258 Questimator questions (each question has one labeler), with most of questions in their areas of expertise. 78% were labeled useful for assessment. Problems with the others included multiple correct answers (36%), irrelevant question topic (22%), obvious answer (21%), all wrong answers (20%) and typos (11%). This kind of labeling has been used in previous work (Section 2). While it gives an idea how effective Questimator is at generating reasonable questions, it does not assess the actual ability of questions to discriminate between different levels of knowledge.

To enhance ecological validity, Questimator assesses actual student performance and compares this to expert-generated quizzes. Note that experimental comparisons with previous work are too difficult to run due to the lack of shared resources. We focused on two key measures:

1. the correlation between a student's performance on a set of Questimator-generated questions, and the same student's performance on a set of human expert generated questions (Section 4.2), and
2. the discriminatory power of individual Questimator-generated questions, in terms of their ability to distinguish between students' with varying topic knowledge states (Section 4.3).

Our motivation for the first objective is that, we would like Questimator to automatically construct a quiz that provides a measure of student knowledge that is similar to an assessment constructed by expert teachers. If Questimator can provide assessments such that a student performance on said assessments correlates highly with the same student's performance on an expert-constructed assessment, that provides encouraging evidence that Questimator is able to capture signals that provide important insight into a student's knowledge.

The second objective stems from wanting a deeper understanding of the quality of the individual questions generated by Questimator. We would like to better quantify how effective different automatically-constructed items are at assessing student knowledge, and how these compare to expert-generated questions. This could also have interesting implications for future work which may generate many questions and then subsample.

4.1 Experimental Setup

Our evaluation was a within-subjects experiment with participants drawn from Amazon Mechanical Turk. In this setup, each learner sees a quiz on a particular topic that is composed of both expert and Questimator generated questions.

Choosing topics We evaluated Questimator on 10 diverse topics. We chose topics for evaluation based on three criteria. First, we chose topics of broad interest because our participants were drawn from Mechanical Turk, which excluded topics like "Reinforcement Learning" that we thought few workers would know about. Second, we chose topics for which we could find existing online quizzes in order to compare to human expert-generated assessments. Finally, we

chose topics that naturally lent themselves to textual questions and answers because Questimator does not yet handle mathematical symbols, images, or video. With these criteria, we chose ten topics across ten disciplines: customer satisfaction (marketing), earthquake (earth science), developmental psychology (psychology), cell (biology), market structure (economics), Vietnam War (history), metaphysics (philosophy), stroke (medicine), waste management (environmental science), elasticity (physics). Quizzes were drawn from MOOCs (Coursera/edX), US university/school board websites, and textbooks by major publishers (*e.g.*, McGraw Hill). We identified two expert-generated quizzes for two topics (Vietnam War and earthquake), and one for the other topics.

Creating quizzes To generate a quiz for a topic, we combined 10 questions from an expert quiz¹ with the 10 Questimator questions. For two topics for which we had two expert quizzes, we generated a quiz which comprised 20 expert questions only, as sampled from both quizzes.

Selecting questions Expert-generated quizzes varied in their length from 10 to 60 questions. To eliminate testing differences across topics, we randomly sampled 10 questions from each expert quiz, after removing questions that relied on numerical calculations and True-False questions wherever possible (if removing these questions resulted in fewer than 10 questions, we retained them). One of the two expert-generated quiz on the "Vietnam war" was drawn from a textbook chapter on the "Vietnam Era", and questions on domestic issues like the Civil Rights Movement were removed before sampling. For each topic, we also generated 10 questions on Questimator.

Participants Participants were recruited from Mechanical Turk. In all, 833 workers participated. All participants were paid \$1 as base for their participation. In addition, they could earn up to \$4 as bonuses based on their test scores.

Experimental procedure Participants were shown one question at a time. To reduce ordering effects, question order was randomized across participants. To reduce response-order biases, the order of answer choices was also randomized. Each question was also augmented with a textbox that asked participants to explain their reasoning (at least 15 characters), a technique that has been previously shown to encourage honest effort [Kittur *et al.*, 2008]. Finally, to discourage participants from using the Internet to search for answers, we monitored the web browser blur event and warned subjects that they would not be allowed to submit an answer if they left the window. Subjects spent an average of 58 seconds on each question. The average length of textbox comments was 66 characters, much longer than the required 15. We received between 78 and 82 completed quizzes for each topic.

4.2 Objective 1: Correlation with Expert Quizzes

We computed a score for both the expert- and Questimator-generated questions for each participant. Each question was weighted equally. The median Pearson correlation between Questimator and expert quiz scores was 0.28 (Table 1).

Of course, even expert-generated quizzes may not correlate highly: two experts may focus on different subtopics within

¹If two expert quizzes were available, we selected one at random

a general topics, or prioritize different forms of knowledge. Therefore, we also evaluated how well scores on two expert quizzes correlate for two topics. The Pearson correlation of the two scores for expert quizzes was 0.430 for “Earthquake” (vs 0.370 for the Questimator-expert) and 0.460 for “Vietnam War” (vs 0.275 for the Questimator-expert).

Scores on Questimator quizzes correlate with expert quizzes to a lower but similar degree as to the two expert quizzes gathered for each category. We believe this means our quiz covers a subset of the whole topic space (also partially covered by expert quiz differently) using reasonable questions. Taken together, it suggests that Questimator scores generally correlate quite well with expert quiz scores.

4.3 Objective 2: Question Discriminative Power

To analyze the discriminative power of questions, we used a very popular approach from psychometrics, Item Response Theory (IRT) [Fox, 2010] which is used to evaluate test items (questions) and analyze test takers. We fit a two-stage IRT model to analyze questions from a quiz under investigation (Questimator quiz) against a reference quiz (expert quiz) when they are mixed together for testing.

The IRT model we use is the unidimensional dichotomous model. For a student i with ability $\theta_i \in \mathbb{R}$, the probability of success on the j th question item is

$$P(Y_j = 1|\theta_i) = \frac{1}{1 + e^{-\alpha_j\theta_i + \beta_j}}, \quad (1)$$

where θ_i is the student ability, α_j and β_j are question parameters. θ_i is a scalar, implying single knowledge ability for the related term affects students’ performance on the quiz. Naturally α_j is also unidimensional. α_j is the question discriminative parameter, and β_j is question difficulty. We selected this simple version of an IRT model to alleviate over-fitting. It can also be easily visualized and interpreted. For the IRT curve of an item (in our case, as question), α specifies the curve steepness (larger α , steeper curve) and β shifts the curve horizontally (larger β , more to the right). We want a steeper curve (in the extreme case, a step function), so there is less randomness in the score given a student’s ability.

More specifically, we considered a Bayesian IRT model with a prior of $\theta_i \sim N(0, 1)$, $\alpha \sim N(\mu_\alpha = 1, \delta_\alpha^2)$, $\beta \sim N(\mu_\beta = 0, \delta_\beta^2)$. μ_α is set to be larger than 0 as we would suppose the questions have some positive discrimination effect. $\mu_\alpha, \mu_\beta, \delta_\alpha$ and δ_β can be set to reflect prior knowledge of the question items. We use MCMC to perform Bayesian inference and estimate the parameters of the model [Fox, 2010].

We fit two IRT models for each topic (*i.e.*, each quiz mixture). We first fit one IRT model with the expert questions. We then treat the estimated student’s ability $\hat{\theta}_{i,exp}$ as the true underlying student abilities of the tested topic. We fix $\theta_i = \hat{\theta}_{i,exp}$ and fit another IRT model with Questimator questions, estimating only α and β for the Questimator questions (Figure 3).

For “Customer satisfaction” (Pearson correlation across overall expert-Questimator question sets: 0.465), all of the questions have positive discrimination ($\alpha > 0$), and a fair amount have relatively large α s (steep sigmoid curve). On the other hand, for “Elasticity (physics)” (correlation: 0.083),

Quiz	Corr	Expert $\bar{\alpha}$	Our $\bar{\alpha}$
Cell (biology)	0.336**	0.834	0.549
Customer satisfaction	0.465**	0.428	0.992
Dev psychology	0.366**	0.992	0.464
Earthquake	0.370**	0.562	0.643
Elasticity (physics)	0.083	0.465	0.267
Market structure	0.282*	0.439	0.616
Metaphysics	0.259*	0.645	0.450
Stroke	0.255*	0.337	0.638
Vietnam War	0.275*	0.838	0.422
Waste management	0.216	0.652	0.548
Average	0.291	0.619	0.559

Table 1: First column: Pearson correlations of quiz scores (**: p-value < 0.01, *: p-value < 0.05). Second and third columns: Mean of the questions’ α s in a quiz.

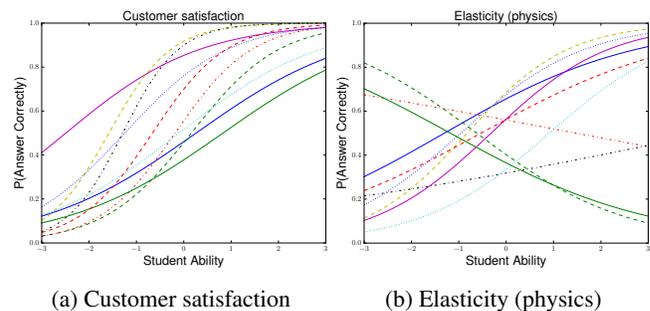


Figure 3: IRT curves of Questimator questions. Each curve corresponds to a question. Larger α , steeper curve.

the questions have lower discriminative power (smaller α s). Few have negative discrimination coefficients ($\alpha < 0$).

Table 1 shows the average α parameters values across the questions for each topic, both for the existing online quizzes and our Questimator-generated quizzes. These results suggest that for many topics Questimator is identifying questions that have positive discriminative power. In the future, it would be interesting to use this approach to automatically refine generated questions or question generation techniques.

5 Conclusion

This paper has introduced Questimator, a system for generating formative, fact-recall questions on arbitrary topics. Our results show that our automatically generated questions are comparable to existing online quizzes on a variety of topics, and that we can generate quizzes for many topics for which no quiz currently exists. Future work will explore (i) improving the quality and the breadth of automatically created assessments, (ii) integrating it into systems for learning through testing, and (iii) utilizing additional corpora.

6 Acknowledgement

The work was supported by a CMU ProSEED grant and a Google focused research award. We also gratefully acknowledge the assistance and/or helpful feedback of our reviewers and participants from MTurk.

References

- [Agarwal and Mannem, 2011] Manish Agarwal and Prashanth Mannem. Automatic gap-fill question generation from text books. In *Proceedings of the 6th Workshop on Innovative Use of NLP for Building Educational Applications*, pages 56–64. ACL, 2011.
- [Bhatia et al., 2013] Arjun Singh Bhatia, Manas Kirti, and Sujjan Kumar Saha. Automatic generation of multiple choice questions using wikipedia. In *Pattern Recognition and Machine Intelligence*, pages 733–738. Springer, 2013.
- [Brandt et al., 2009] Joel Brandt, Philip J Guo, Joel Lewenstein, Mira Dontcheva, and Scott R Klemmer. Two studies of opportunistic programming: interleaving web foraging, learning, and writing code. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1589–1598. ACM, 2009.
- [Brown et al., 2005] Jonathan C Brown, Gwen A Frishkoff, and Maxine Eskenazi. Automatic question generation for vocabulary assessment. In *Proceedings of the conference on HLT and EMNLP*, pages 819–826. ACL, 2005.
- [Cai et al., 2015] Carrie J Cai, Philip J Guo, James R Glass, and Robert C Miller. Wait-learning: Leveraging wait time for second language education. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 3701–3710. ACM, 2015.
- [Caputo and Dunning, 2005] Deanna Caputo and David Dunning. What you dont know: The role played by errors of omission in imperfect self-assessments. *Journal of Experimental Social Psychology*, 41(5):488–505, 2005.
- [Christoforaki and Ipeirotis, 2015] Maria Christoforaki and Panagiotis G Ipeirotis. A system for scalable and reliable technical-skill testing in online labor markets. *Computer Networks*, 2015.
- [Elwood and Klenowski, 2002] Jannette Elwood and Val Klenowski. Creating communities of shared practice: the challenges of assessment use in learning and teaching. *Assessment & Evaluation in Higher Education*, 27(3):243–256, 2002.
- [Fox, 2010] Jean-Paul Fox. *Bayesian item response modeling: Theory and applications*. Springer Science & Business Media, 2010.
- [Gates et al., 2011] Donna Marie Gates, Greg Aist, Jack Mostow, Margaret McKeown, and Juliet Bey. How to generate cloze questions from definitions: A syntactic approach. In *2011 AAAI Fall Symposium Series*, 2011.
- [Heilman and Smith, 2010] Michael Heilman and Noah A Smith. Good question! statistical ranking for question generation. In *NAACL/HLT 2010*. ACL, 2010.
- [Horwitz and Cope, 1986] Horwitz M.B. Horwitz, E.K. and J. Cope. Foreign language classroom anxiety. *The Modern language journal*, 70(2):125–132, 1986.
- [Hoshino and Nakagawa, 2005] Ayako Hoshino and Hiroshi Nakagawa. A real-time multiple-choice question generation for language testing: A preliminary study. In *EdApp-NLP*, pages 17–20, Stroudsburg, PA, USA, 2005. ACL.
- [Kiros et al., 2015] Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Skip-thought vectors. *arXiv preprint arXiv:1506.06726*, 2015.
- [Kittur et al., 2008] Aniket Kittur, Ed H Chi, and Bongwon Suh. Crowdsourcing user studies with mechanical turk. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 453–456. ACM, 2008.
- [Kittur et al., 2009] Aniket Kittur, Ed H Chi, and Bongwon Suh. What’s in wikipedia?: mapping topics and conflict using socially annotated category structure. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 1509–1512. ACM, 2009.
- [Kulkarni et al., 2015] Chinmay Kulkarni, Julia Cambre, Yasmine Kotturi, Michael S Bernstein, and Scott R Klemmer. Talkabout: Making distance matter with small groups in massive classes. In *Proceedings of the 18th ACM CSCW*, pages 1116–1128. ACM, 2015.
- [Lin et al., 2007] Yi-Chien Lin, Li-Chun Sung, and Meng Chang Chen. An automatic multiple-choice question generation scheme for english adjective understanding. In *Workshop on Modeling, Management and Generation of Problems/Questions in eLearning, the 15th ICCE*, pages 137–142, 2007.
- [Mikolov et al., 2013] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [Miller, 1995] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [Mitkov et al., 2006] Ruslan Mitkov, Le An Ha, and Niki-foros Karamanis. A computer-aided environment for generating multiple-choice test items. *Natural Language Engineering*, 12(02):177–194, 2006.
- [Mostow and Jang, 2012] Jack Mostow and Hyeju Jang. Generating diagnostic multiple choice comprehension cloze questions. In *Proceedings of the Seventh Workshop on Building Educational Applications Using NLP*, pages 136–146. ACL, 2012.
- [Rohde, 2004] Douglas LT Rohde. Tgrep2 user manual, 2004.
- [Sadowski et al., 2015] Caitlin Sadowski, Kathryn T Stolee, and Sebastian Elbaum. How developers search for code: a case study. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*, pages 191–201. ACM, 2015.
- [Smith et al., 2010] Simon Smith, PVS Avinesh, and Adam Kilgarrieff. Gap-fill tests for language learners: Corpus-driven item generation. In *Proceedings of ICON*, 2010.
- [Wang et al., 2008] Weiming Wang, Tianyong Hao, and Wenyin Liu. Automatic question generation for learning evaluation in medicine. In *Advances in Web Based Learning-ICWL 2007*, pages 242–251. Springer, 2008.